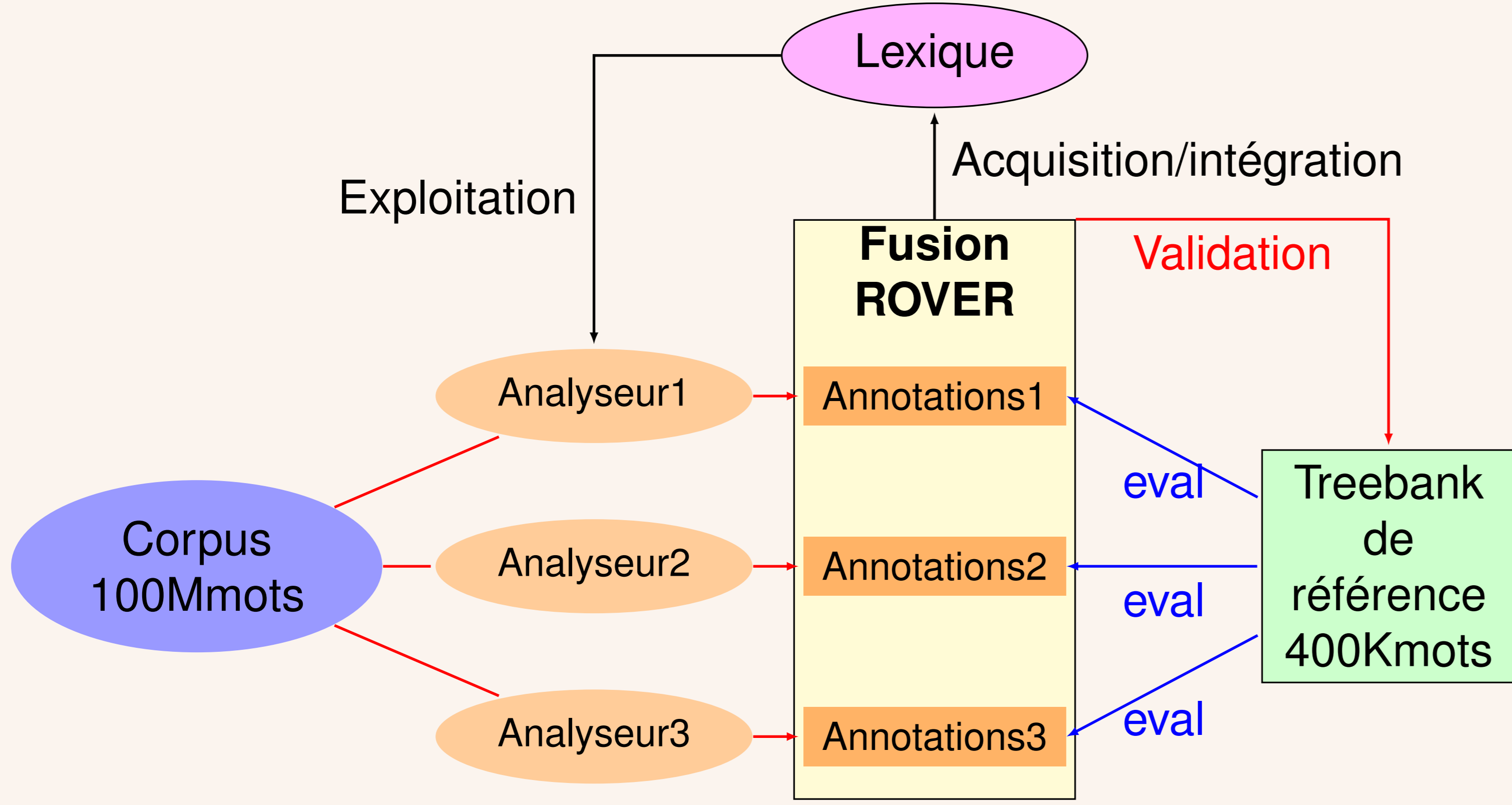


Établir un cercle vertueux entre outils et ressources linguistiques



Exploiter de très grands corpus

Les corpus annotés syntaxiquement comme le **Penn Tree Bank** sont très importants pour le **Traitement Automatique des Langues [TAL]** mais **rares** et **coûteux** à développer.

Par contre, il est maintenant possible d'accéder à de très grandes quantités de textes électroniques en français:

Corpus	Taille	Nature
Corpus EASy	1Mmots	multi-styles
Wikipedia Fr	~ 86Mmots	encyclopédique collaboratif libre
Wikisources	~ 80Mmots	littéraire libre
Monde Diplomatique	18Mmots	journalistique
FRANTEXT	20Mmots	littéraire libre
Europarl	28Mmots	débat Parlement européen
JRC-Acquis	39Mmots	juridique européen
Corpus Ester	1Mmots	oral transcrit
Total (actuel)	> 270 Mmots	

Profiter d'une dizaine d'analyseurs

Une occasion unique, source de diversité (formalismes, technologies, ...)

Analyseur	Origine	Nature
FRMG	INRIA	TIG/TAG+DYALOG
SxLFG	INRIA	LFG+SYNTAX
LLP2	LORIA	TAG
LIMA	CEA-LIST	
TAGPARSER	TAGMATICA	Induction + règles
SYNTAX	ERSS	Système de règles
GP1 & GP2	LPL	Grammaires de Propriétés
CORDIAL	SYNAPSE	
SYGMART	LIRMM	
XIP	XRCE	Cascade de règles

Évaluer et combiner les annotations

- ▶ Expertise de Technolangue **IVALDA/EASy** sur l'évaluation des analyseurs du français:
 - [Eval1] Campagne Novembre 2007
 - [Eval2] Campagne fin 2009
- ▶ Format EASy: chunks + dépendances
 - Ajout des groupes récursifs (constituance)
 - Vers le respect des **standards** ISO TC37SC4.
- ▶ **ROVER**: Fusion des jeux d'annotations paramétrée par les évaluations + *feedback*
- ▶ Validation manuelle d'un sous-corpus
- ▶ ⇒ Améliorer les analyseurs !

Acquérir et intégrer des connaissances

- ▶ Exploiter les annotations ROVER pour de l'acquisition de **connaissances lexicales**:
 - information de valence (verbes, ...)
 - classes d'alternation verbales (Levin)
 - probabilités de désambiguïation catégories, constructions, ...
 - restrictions de sélection
 - classes sémantiques (Hyp. distributionnelle Harris)
 - morphologie dérivationnelle avec transfert des informations syntaxiques
- ▶ Les valider et intégrer dans un lexique
- ▶ Les exploiter dans certains analyseurs ⇒ Améliorer les analyseurs !

Calendrier indicatif

Corpus	Format	Analyse Corpus	Distrib
	Eval1	Fusion Annotations	
		Validation Treebank	
		Acquisition	
		Intégration	
	Eval2		

Infrastructure pour passer à l'échelle

Besoin d'une infrastructure solide pour gérer des masses de données

- ▶ Utilisation de *fermes* de machines pour analyser de gros corpus
- ▶ **EASYREF**, un prototype WEB de gestion collaborative d'annotations *visualisation, recherche, édition, versionning, comparaison, dépôt, fusion, ...*

▶ Rapports pour oral_delic_1:E96

Actions	Rid	Bid	Corpus:Phrase:Rev	Auteur	Date	S C T
Del Edit Next	645	645	oral_delic_1 : E96 : r000	gil	2007-08-09 13:29:30	o
Del Edit Next	979	645	oral_delic_1 : E96 : r000	nora	2007-10-15 12:30:13	x

New

▶ Erreurs potentielles détectées pour oral_delic_1:E96

▶ Corrections pour oral_delic_1:E96

Rev0001 created relation E96R4 with type 'COMP' and roles 'complementeur' => 'E96F6' 'verbe' => 'E96F10' -- nora -- lun

Annotations pour oral_delic_1:E96

S	D	Enoncé oral_delic_1 E96 -- Rev0001 -- Analyse complète FRMG											
		NV 1	GP 2		NV 3	GR 4	NV 5	GN 6					
		NV 1	GP 2	GN 3	NV 4	GR 5	NV 6	GN 7					
		donc	c'	est	pour	ça	qu'	elle	a	pas	eu	son	C.A.P
		1	2	3	4	5	6	7	8	9	10	11	12
		COORD			MOD-N								
		SUJ-V		COD-V									
		ATB-SO				COMP							
		CPL-V			SUJ-V		MOD-V						
				AUX-V		MOD-V		COD-V					

Retombées attendues

- ▶ Rendre librement disponible des ressources linguistiques de qualité pour le français:
 - des corpus annotés, dont
 - ▶ le **ROVER** (> 100 Mmots)
 - ▶ un **treebank** vérifié manuellement (400Kmots)
 - des ressources lexicales
- ▶ Améliorer:
 - les analyseurs syntaxiques robustes efficaces du français
 - les méthodologies d'évaluation et de fusion des annotations
 - les méthodologies d'acquisition par bootstrap
- ▶ Permettre des applications par analyse syntaxique de corpus:
 - Acquisition de connaissances (extraction d'ontologies, ...)
 - Extraction d'informations (fouille de textes, veille, question-réponse, ...)